code for checking the Quality of Wine

```python
In [3]: import numpy as np # linear algebra
        import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
        import plotly.express as px
        import seaborn as sns
        import matplotlib.pyplot as plt
        import plotly.graph_objects as go
        from plotly.subplots import make_subplots

        from xgboost import XGBClassifier,XGBRegressor
```

```python
In [12]: import os
         for dirname, _, filenames in os.walk("C:\\Users\\amith\\Downloads\\winequality-red.csv"):
             for filename in filenames:
                 print(os.path.join(dirname, filename))
```

```python
In [10]: # Importing the necessary libraries
         import pandas as pd
         import numpy as np
         import matplotlib.pyplot as plt
         import seaborn as sns
         from sklearn.model_selection import train_test_split
         from sklearn.linear_model import LinearRegression
         from sklearn.metrics import mean_squared_error, r2_score
```

```python
In [11]: # Load the dataset (winequality-red.csv or winequality-white.csv)
         url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv'
         data = pd.read_csv(url, sep=';')
```

```python
In [12]: # Display the first 5 rows of the dataset
         print(data.head())
```

```
   fixed acidity  volatile acidity  citric acid  residual sugar  chlorides  \
0            7.4              0.70         0.00             1.9      0.076
1            7.8              0.88         0.00             2.6      0.098
2            7.8              0.76         0.04             2.3      0.092
3           11.2              0.28         0.56             1.9      0.075
4            7.4              0.70         0.00             1.9      0.076

   free sulfur dioxide  total sulfur dioxide  density    pH  sulphates  \
0                 11.0                  34.0   0.9978  3.51       0.56
1                 25.0                  67.0   0.9968  3.20       0.68
2                 15.0                  54.0   0.9970  3.26       0.65
3                 17.0                  60.0   0.9980  3.16       0.58
4                 11.0                  34.0   0.9978  3.51       0.56

   alcohol  quality
0      9.4        5
1      9.8        5
2      9.8        5
3      9.8        6
4      9.4        5
```

```python
In [13]: # Data overview
         print(data.info())
         print(data.describe())

         # Check for missing values
         print(data.isnull().sum())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1599 entries, 0 to 1598
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype
---  ------                --------------  -----
 0   fixed acidity         1599 non-null   float64
 1   volatile acidity      1599 non-null   float64
 2   citric acid           1599 non-null   float64
 3   residual sugar        1599 non-null   float64
 4   chlorides             1599 non-null   float64
 5   free sulfur dioxide   1599 non-null   float64
 6   total sulfur dioxide  1599 non-null   float64
 7   density               1599 non-null   float64
 8   pH                    1599 non-null   float64
 9   sulphates             1599 non-null   float64
 10  alcohol               1599 non-null   float64
 11  quality               1599 non-null   int64
dtypes: float64(11), int64(1)
memory usage: 150.0 KB
None
       fixed acidity  volatile acidity  citric acid  residual sugar  \
count    1599.000000       1599.000000  1599.000000     1599.000000
mean        8.319637          0.527821     0.270976        2.538806
std         1.741096          0.179060     0.194801        1.409928
min         4.600000          0.120000     0.000000        0.900000
25%         7.100000          0.390000     0.090000        1.900000
50%         7.900000          0.520000     0.260000        2.200000
75%         9.200000          0.640000     0.420000        2.600000
max        15.900000          1.580000     1.000000       15.500000

          chlorides  free sulfur dioxide  total sulfur dioxide      density  \
count  1599.000000          1599.000000           1599.000000  1599.000000
mean      0.087467            15.874922             46.467792     0.996747
std       0.047065            10.460157             32.895324     0.001887
min       0.012000             1.000000              6.000000     0.990070
25%       0.070000             7.000000             22.000000     0.995600
50%       0.079000            14.000000             38.000000     0.996750
75%       0.090000            21.000000             62.000000     0.997835
max       0.611000            72.000000            289.000000     1.003690

                pH    sulphates      alcohol      quality
count  1599.000000  1599.000000  1599.000000  1599.000000
mean      3.311113     0.658149    10.422983     5.636023
std       0.154386     0.169507     1.065668     0.807569
min       2.740000     0.330000     8.400000     3.000000
25%       3.210000     0.550000     9.500000     5.000000
50%       3.310000     0.620000    10.200000     6.000000
75%       3.400000     0.730000    11.100000     6.000000
max       4.010000     2.000000    14.900000     8.000000
fixed acidity           0
volatile acidity        0
citric acid             0
residual sugar          0
chlorides               0
free sulfur dioxide     0
total sulfur dioxide    0
density                 0
pH                      0
sulphates               0
alcohol                 0
quality                 0
dtype: int64
```

```python
In [14]: # Split data into features (X) and target (y)
         X = data.drop('quality', axis=1)
         y = data['quality']
```

```python
In [18]: # Split the dataset into training and testing sets (80% train, 20% test)
         X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```python
In [19]: from sklearn.ensemble import RandomForestRegressor
         from sklearn.model_selection import cross_val_score

         # Initialize Random Forest model
         rf_model = RandomForestRegressor(n_estimators=100, random_state=42)

         # Train the model
         rf_model.fit(X_train, y_train)

         # Predict on test set
         y_pred_rf = rf_model.predict(X_test)
```

```python
In [20]: # Evaluate the model
         rf_r2 = r2_score(y_test, y_pred_rf)
         rf_mse = mean_squared_error(y_test, y_pred_rf)
         rf_rmse = np.sqrt(rf_mse)

         print(f'Random Forest R² Score: {rf_r2}')
         print(f'Random Forest MSE: {rf_mse}')
         print(f'Random Forest RMSE: {rf_rmse}')
```

```
Random Forest R² Score: 0.5390429623873638
Random Forest MSE: 0.3012381249999997
Random Forest RMSE: 0.5488516420673258
```

```python
In [23]: # Feature importance analysis
         coefficients = pd.DataFrame(model.coef_, X.columns, columns=['Coefficient'])
         print(coefficients)
```

```
                      Coefficient
fixed acidity            0.022593
volatile acidity        -1.042144
citric acid             -0.103974
residual sugar           0.013310
chlorides               -1.764144
free sulfur dioxide      0.004627
total sulfur dioxide    -0.003248
density                -17.239629
pH                      -0.342743
sulphates                0.876382
alcohol                  0.271549
```

```python
In [24]: # Assign values for the features (as a single data point)
         # Example: Create a dictionary with custom feature values
         new_data = {
             'fixed acidity': 7.4,
             'volatile acidity': 0.7,
             'citric acid': 0.0,
             'residual sugar': 1.9,
             'chlorides': 0.076,
             'free sulfur dioxide': 11.0,
             'total sulfur dioxide': 34.0,
             'density': 0.9978,
             'pH': 3.51,
             'sulphates': 0.56,
             'alcohol': 9.4
         }

         # Convert the dictionary to a DataFrame (single row of feature values)
         new_data_df = pd.DataFrame([new_data])

         # Display the new data point
         print("New data point:")
         print(new_data_df)

         # Use the trained model to predict wine quality for the new data point
         predicted_quality = model.predict(new_data_df)

         # Display the predicted quality
         print(f'\nPredicted Wine Quality: {predicted_quality[0]}')
```

```
New data point:
   fixed acidity  volatile acidity  citric acid  residual sugar  chlorides  \
0            7.4               0.7          0.0             1.9      0.076

   free sulfur dioxide  total sulfur dioxide  density    pH  sulphates  \
0                 11.0                  34.0   0.9978  3.51       0.56
```

```
   alcohol
0     9.4
```

Predicted Wine Quality: 5.040537396744657